

We don't need no bounding-boxes: Training object class detectors using only human verification

Dim P. Papadopoulos

dim.papadopoulos@ed.ac.uk

Jasper R. R. Uijlings

jrr.uijlings@ed.ac.uk

Frank Keller

keller@inf.ed.ac.uk

Vittorio Ferrari

vferrari@inf.ed.ac.uk

University of Edinburgh

Abstract

Training object class detectors typically requires a large set of images in which objects are annotated by bounding-boxes. However, manually drawing bounding-boxes is very time consuming. We propose a new scheme for training object detectors which only requires annotators to verify bounding-boxes produced automatically by the learning algorithm. Our scheme iterates between re-training the detector, re-localizing objects in the training images, and human verification. We use the verification signal both to improve re-training and to reduce the search space for re-localisation, which makes these steps different to what is normally done in a weakly supervised setting. Extensive experiments on PASCAL VOC 2007 show that (1) using human verification to update detectors and reduce the search space leads to the rapid production of high-quality bounding-box annotations; (2) our scheme delivers detectors performing almost as good as those trained in a fully supervised setting, without ever drawing any bounding-box; (3) as the verification task is very quick, our scheme substantially reduces total annotation time by a factor $6 \times -9 \times$.

1. Introduction

Object class detection is a central problem in computer vision. Training a detector typically requires a large set of images in which objects have been manually annotated with bounding-boxes [10, 16, 18, 19, 20, 34, 53, 57, 62]. Bounding-box annotation is tedious, time consuming and expensive. For instance, annotating ILSVRC, currently the most popular object class detection dataset, required 42s per bounding-box by crowd-sourcing on Mechanical Turk [40] using a technique specifically developed for efficient bounding-box annotation [50].

In order to reduce the cost of bounding-box annotation, researchers have mainly focused on two strategies. The first is learning in the weakly supervised setting, i.e. given only labels indicating which object classes are present in an image. While this setting is much cheaper, it produces lower quality detectors, typically performing only about

half as well as when trained from bounding-boxes [4, 5, 9, 12, 42, 47, 48, 49, 60]. The second strategy is active learning, where the computer requests human annotators to draw bounding-boxes on a subset of images actively selected by the learner itself. This strategy can produce high quality detectors, but it still requires humans to draw a lot of bounding-boxes in order to get there, leading to limited gains in terms of total annotation time [56, 63].

In this paper we propose a new scheme for learning object detectors which only requires humans to *verify* bounding-boxes produced automatically by the learning algorithm: the annotator merely needs to decide whether a bounding-box is correct or not. Crucially, answering this verification question takes much less time than actually drawing the bounding-box.

Given a set of training images with image-level labels, our scheme iteratively alternates between updating object detectors, re-localizing objects in the training images, and querying humans for verification. At each iteration we use the verification signal in two ways. First, we update the object class detector using only positively verified bounding-boxes. This makes it stronger than when using all detected bounding-boxes, as it is commonly done in the weakly supervised setting, because typically many of them are incorrect. Moreover, once the object location in an image has been positively verified, it can be fixed and removed from consideration in subsequent iterations. Second, we observe that bounding-boxes judged as incorrect still provide valuable information about where the object is not. Building on this observation, we use the negatively verified bounding-boxes to reduce the search space of possible object locations in subsequent iterations. Both these points help to rapidly find more objects in the remaining images. This results in a framework for training object detectors which minimizes human annotation effort and eliminates the need to draw any bounding-box.

In extensive experiments on the popular PASCAL VOC 2007 dataset [16] with both simulated and actual annotators, we show that: (1) using human verification to update detectors and reduce the search space leads to rapid pro-

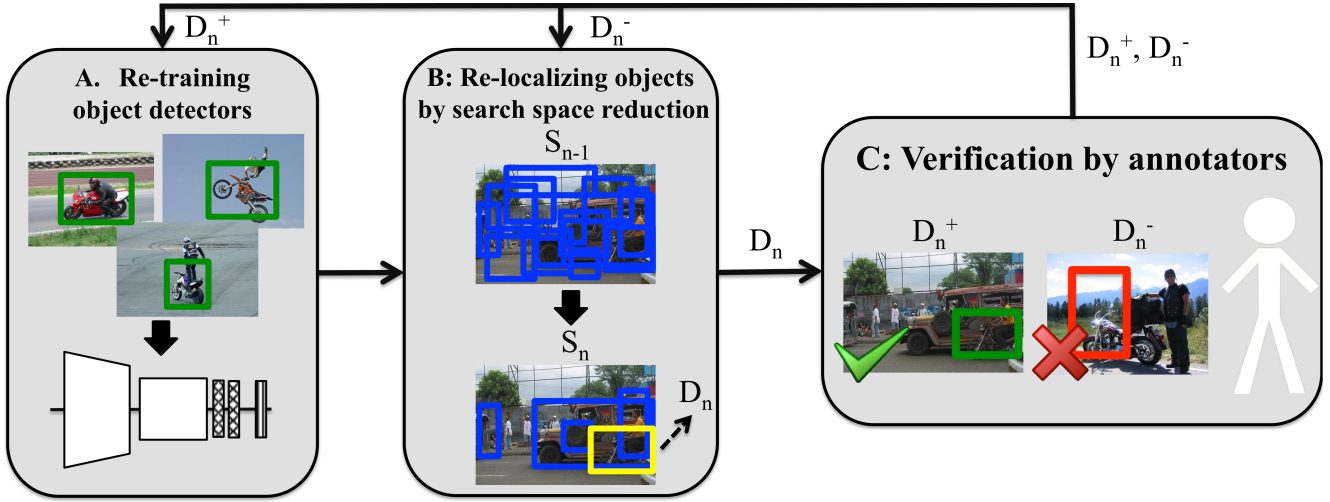


Figure 1: Our framework iterates between (A) re-training object detectors, (B) re-localising objects, and (C) querying annotators for verification. The verification signal resulting from (C) is used in both (A) and (B).

duction of high-quality bounding-box annotations; (2) our scheme delivers object class detectors performing almost as good as those trained in a fully supervised setting, *without ever drawing any bounding-box*; (3) as the verification task is very quick, our scheme substantially reduces total annotation time by a factor $6 \times 9 \times$.

2. Related Work

Weakly-supervised object localization (WSOL). Many previous techniques [4, 5, 8, 12, 42, 47, 48, 49, 60] try to learn object class detectors in the weakly supervised setting, i.e. given training images known to contain instances of a certain object class, but not their location. The task is both to localize the objects in the training images and to learn an object detector for localizing instances in new images.

Recent work on WSOL [4, 5, 8, 48, 49, 60] has shown remarkable progress, mainly because of the use of Convolutional Neural Nets (CNN [20, 28]), which greatly improve visual recognition tasks. However, learning a detector without location annotation is difficult and performance is still quite low compared to fully supervised methods (typically about half the mAP of the same detection model trained on the same images but with manual bounding-box annotation [4, 5, 8, 12, 42, 47, 48, 49, 60]).

Often WSOL is conceptualised as Multiple Instance Learning (MIL) [4, 8, 12, 13, 44, 47, 48, 49]. Images are treated as bags of windows (instances). A negative image contains only negative instances. A positive image contains at least one positive instance, mixed in with a majority of negative ones. The goal is to find the true positive instances from which to learn a window classifier for the object class. This typically happens by iteratively alternating between (A) re-training the object detector given the current selection of positive instances and (B) re-localising instances in the positive images using the current object detector. Our

proposed scheme also contains these steps. However, we introduce a human verification step, whose signal is fed into both (A) and (B), which fundamentally alters these steps. The resulting framework leads to substantially better object detectors with modest additional annotation effort (sec. 5.3).

Humans in the loop. Human-machine collaboration approaches have been successfully used in tasks that are currently too difficult to be solved by computer vision alone, such as fine-grained visual recognition [7, 11, 58, 59], semi-supervised clustering [30], attribute-based image classification [6, 36, 37]. These works combine the responses of pre-trained computer vision models on a new test image with human input to fully solve the task. In the domain of object detection, Russakovsky et al. [41] propose such a scheme to fully detect all objects in images of complex scenes. Importantly, their object detectors are pre-trained on bounding-boxes from the large training set of ILSVRC 2014 [40], as their goal is not to make an efficient training scheme.

Active learning. Active learning schemes iteratively train models while requesting humans to annotate a subset of the data points actively selected by the learner as being the most informative. Previous active learning work has mainly focussed on image classification [25, 26, 27, 39], and free-form region labelling [45, 54, 55].

A few works have proposed active learning schemes specifically for training object class detectors [56, 63]. Vijayanarasimhan and Grauman [56] propose an approach where the training images do not come from a predefined dataset but are crawled from the web. Here annotators are asked to draw many bounding-boxes around the target objects (about one third of the training images [56]). Yao et al. [63] propose to manually correct bounding-boxes detected in video. While both [56, 63] produce high quality detectors, they achieve only moderate gains in annotation

time, because drawing or correcting bounding-boxes is expensive. In contrast, our scheme only asks annotators to verify bounding-boxes, never to draw. This leads to more substantial reductions in annotation time.

Other ways to reduce annotation effort. A few authors tried to learn object detectors from videos, where the spatio-temporal coherence of the video frames facilitates object localization [32, 38, 52]. An alternative is transfer learning, where learning a model for a new class is helped by labeled examples of related classes [2, 17, 21, 29, 31]. Hoffman et al. [24] proposed an algorithm that transforms an image classifier to an object detector without bounding-box annotated data using domain adaptation. Other types of data, such as text from web pages or newspapers [3, 15, 22, 33] or eye-tracking data [35], have also been used as a weak annotation signal to train object detectors.

3. Method

In this paper we are given a training set with image-level labels. Our goal is to obtain object instances annotated by bounding-boxes and to train good object detectors while minimizing human annotation effort. We therefore propose a framework where annotators only need to *verify* bounding-boxes automatically produced by our scheme.

Our framework iteratively alternates between (A) re-training object detectors, (B) re-localizing objects in the training images, and (C) querying annotators for verification (fig. 1). Importantly, we use verification signals to help both re-training and re-localisation.

More formally, let I_n be the set of images for which we do not have positively verified bounding-boxes at iteration n yet. Let S_n be the corresponding set of possible object locations. Initially, I_0 is the complete training set and S_0 is a complete set of object proposals [1, 14, 53] extracted from these images (we use EdgeBoxes [14]). To facilitate exposition, we describe our framework starting from the verification step (C, sec 3.1). At iteration n we have a set of automatically detected bounding-boxes D_n which are given to annotators to be verified. Detections which are judged to be correct $D_n^+ \subseteq D_n$ are used for re-training the object detectors (A) in the next iteration (sec. 3.2). The verification signal is also used to reduce the search space S_{n+1} for re-localisation (B, sec. 3.3). We describe our main three steps below. We defer to Sec. 4 a description of the object detection model we use, and of how to automatically obtain initial detections D_0 to start the process.

3.1. Verification by annotators

In this phase, we ask annotators to verify the automatically generated detections D_n at iteration n . For this we explore two strategies (fig. 2): simple yes/no verification, and more elaborate verification in which annotators are asked to categorize the type of error.



Figure 2: Our two verification strategies for some images of the dog class. Yes/No verification (left): verify a detection as either correct (Yes) or incorrect (No). YPCMM verification (right): label a detection as Yes, Part, Container, Mixed or Missed.

Yes/No Verification. In this task the annotators are shown a detection l_d and a class label. They are instructed to respond Yes if the detection correctly localizes an object of that class, and No otherwise. This splits the set of object detections D_n into D_n^+ and D_n^- . We define “correct localization” based on the standard PASCAL Intersection-over-Union criterion [16] (IoU). Let l_d be the detected object bounding-box and l_{gt} be the actual object bounding-box (which is not given to the annotator). Let $\text{IoU}(l_a, l_b) = |l_a \cap l_b| / |l_a \cup l_b|$, where $|\cdot|$ denotes area. If $\text{IoU}(l_{gt}, l_d) \geq 0.5$, the detected bounding-box should be considered correct and the annotator should answer Yes. Intuitively, this Yes/No verification is a relatively simple task which should translate into fast annotation times.

Yes/Part/Container/Mixed/Missed Verification. In this task, we asked the annotators to label an object detection l_d as Yes (correct), Part, Container, Mixed, or Missed. Yes is defined as above ($\text{IoU}(l_{gt}, l_d) \geq 0.5$). For incorrect detections the annotators are asked to diagnose the error as either *Part* if it contains part of the target object and no background; *Container* if it contains the whole object and some background; *Mixed* if it contains part of the object and some background; *Missed* if the object was completely missed. This verification step splits D_n into D_n^+ and $D_n^{\text{ypcmm}-}$. Intuitively, determining the type of error is more difficult leading to longer annotation times, but also brings more information that we can use in the next steps.

3.2. Re-training object detectors

In this step we re-train object detectors. After the verification step we know that D_n^+ contains well localized object instances, while D_n^- or $D_n^{\text{ypcmm}-}$ do not. Hence we train using only bounding-boxes $D_1^+ \cup \dots \cup D_n^+$ that have been positively verified in past iterations. To obtain background training samples, we sample proposals which have an IoU in range $[0 - 0.5)$ with positively verified bounding-boxes.

Note how it is common in WSOL works [4, 5, 8, 12, 42, 47, 48, 49, 60] to also have a re-training step. However, they typically use all detected bounding-boxes D_n . Since in WSOL generally less than half of them are correct, this leads to rather weak detectors. In contrast, our verification

step enables us to train purely from correct bounding-boxes, resulting in stronger, more reliable object detectors.

3.3. Re-localizing objects by search space reduction

In this step we re-localize objects in the training images. For each image, we apply the current object detector to score the object proposals in it, and select the proposal with the highest score as the new detection for that image. Importantly, we do not evaluate *all* proposals S_0 , but instead use the verification signal to reduce the search space by removing proposals.

Positively verified detections D_n^+ are correct by definition and therefore their images need not be considered in subsequent iterations, neither in the re-localization step nor in the verification step. For negatively verified detections we reduce the search space depending on the verification strategy, as described below.

Yes/No Verification. In the case where the annotator judges a detection as incorrect (D_n^-), we can simply eliminate its proposal from the search space. This results in the updated search space S_{n+1} , where one proposal has been removed from each image with an incorrect detection.

However, we might make a better use of the negative verification signal. Since an incorrect detection has an IoU < 0.5 with the true bounding-box, we can eliminate all proposals with an IoU ≥ 0.5 with it. This is a more aggressive reduction of the search space. While it may remove some proposals which are correct according to the IoU criterion, it will not remove the best possible proposal. Importantly, this strategy eliminates those areas of the search space that matter: high scoring locations which are unlikely to contain the object. In Sec. 5.2 we investigate which way of using negatively verified detection performs better in practice.

Yes/Part/Container/Mixed/Missed Verification. In the case where annotators categorize incorrect detections as Part/Container/Mixed/Missed, we can use the type of error to get an even greater reduction of the search space. Depending on the type of error we eliminate different proposals (fig. 3): *Part*: eliminate all proposals which do not contain the detection; *Container*: eliminate all proposals which are not inside the detection; *Mixed*: eliminate all proposals which are not inside the detection, or do not contain it, or have zero IoU with it, or have IoU ≥ 0.5 with it; *Missed*: eliminate all proposals which have non-zero IoU with the detection.

To precisely determine what is “inside” and “contained”, we introduce the Intersection-over-A measure: $\text{IoA}(l_a, l_b) = |l_a \cap l_b| / |l_a|$. Note that $\text{IoA}(l_{gt}, l_d) = 1$ if the detection l_d contains the true object bounding-box l_{gt} , whereas $\text{IoA}(l_d, l_{gt}) = 1$ if l_d covers a part of l_{gt} . In practice, if l_d is judged to be a *Part* by the annotator, we eliminate all proposals l_s with $\text{IoA}(l_d, l_s) \leq 0.9$. Similarly, if l_d is judged to be a *Container*, we eliminate all proposals l_s

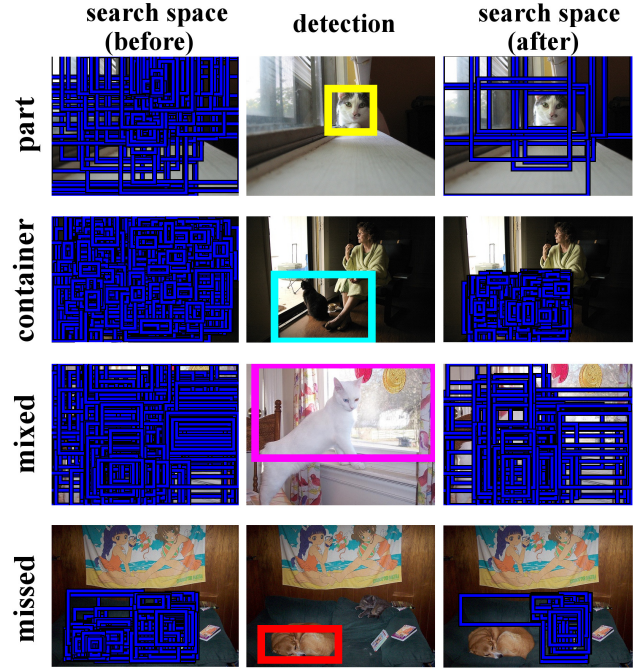


Figure 3: Visualisation of **search space** reduction induced by YPCM verification on some images of the cat class (**part**, **container**, **mixed**, and **missed**). In the last row, the search space reduction steers the re-localization process towards the small cat on the right of the image and away from the dog on the left.

with $\text{IoA}(l_s, l_d) \leq 0.9$. We keep the tolerance threshold 0.9 fixed in all experiments.

Note how in WSOL there is also a re-localisation step. However, because there is no verification signal there is also no search space reduction: each iteration needs to consider the complete set S_0 of proposals. In contrast, in our work the search space reduction greatly facilitates re-localization.

4. Implementation details

We summarize here two existing state-of-the-art components that we use in our framework: the object detection model, and a WSOL algorithm which we use to obtain initial object detections D_0 .

4.1. Object class detector

As object detector we use Fast R-CNN [19], which combines object proposals [1, 14, 53] with CNNs [23, 28, 46]. Instead of Selective Search [53] we use EdgeBoxes [14] which gives us an “objectness” measure [1] which we use in the initialization phase described below. For simplicity of implementation, for the re-training step (sec. 3.2) we omit bounding-box regression, so that the set of object proposals stays fixed throughout all iterations. For evaluation on the test set, we then train detectors with bounding-box regression. In most of our experiments, we use AlexNet [28] as the underlying CNN architecture.

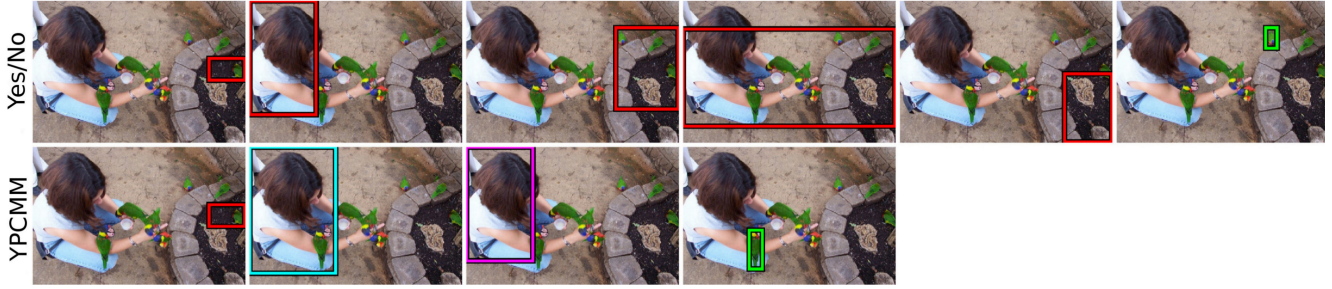


Figure 4: Comparing the search process of Yes/No verification with Yes/Part/Container/Mixed/Missed on an image of the bird class. The extra signal for YPCMM allows for a more targeted search, resulting in fewer verification steps to find the object.

4.2. Initialization by Multiple Instance Learning

We perform multiple instance learning (MIL) for weakly supervised object localisation [4, 8, 48] to obtain the initial set of detections D_0 . We start with the training images I_0 and the set of object proposals S_0 extracted using EdgeBoxes [14]. Following [4, 20, 48, 49, 60] we extract CNN features on top of which we train an SVM. We iterate between (A) re-training object detectors and (B) re-localizing objects in the training images. We stop when two subsequent re-localization steps yield the same detections, which typically happens within 10 iterations. These detections become D_0 . In the very first iteration, we train the classifier using complete images as positive training examples [8, 42].

We apply two improvements to the standard MIL framework. First, in high dimensional feature space the discriminative SVM classifier can relatively easily separate any positive examples from negative examples, which means that most positive examples are far from the decision hyperplane. Hence the same positive training examples used for re-training (A) are often re-localised in (B), leading to premature locked-in behaviour. To prevent this Cinbis et al. [8, 9] introduced multi-fold MIL: similar to cross-validation, the dataset is split into 10 subsets, where the re-localisation on each subset is done using detectors trained on the union of all other subsets. Second, like in [9, 12, 21, 38, 43, 47, 44, 51, 61], we combine the object detector score with a general measure of “objectness” [1], which measures how likely it is that a proposal tightly encloses an object of any class (e.g. bird, car, sheep), as opposed to background (e.g. sky, water, grass). In this paper we use the recent objectness measure of [14].

5. Experimental Results

5.1. Dataset and evaluation protocol

PASCAL VOC 2007. We perform experiments on PASCAL VOC 2007 [16], which consists of 20 classes. The trainval set contains 5011 images, while the test set contains 4952 images. We use the trainval set with accompanying image-level labels to train object detectors, and measure their performance on the test set. Following the common protocol for WSOL experiments [8, 9, 12, 42, 60],

we exclude trainval images that contain only difficult and truncated instances, ending up with 3550 images. In sections 5.2, 5.3 we carry out a detailed analysis of our system in these settings, using AlexeNet as CNN architecture [28]. For completeness, in section 5.4 we also present results when using the complete trainval set and the deeper VGG16 as CNN architecture [46].

Evaluation. Given a training set with image-level labels, our goal is to localize the object instances in this set and to train good object detectors, while minimizing human annotation effort. We evaluate this by exploring the trade-off between localization performance and quality of the object detectors versus required annotation effort. We quantify localization performance in the training set with the Correct Localization (CorLoc) measure [4, 5, 8, 9, 12, 42, 47, 60]. CorLoc is the percentage of images in which the bounding-box returned by the algorithm correctly localizes an object of the target class (i.e., $\text{IoU} \geq 0.5$).

We quantify object detection performance on the test set using mean average precision (mAP), as standard in PASCAL VOC 07. We quantify annotation effort both in terms of the number of verifications and in terms of actual human time measurements.

As most previous WSOL methods [4, 5, 8, 9, 12, 42, 47, 48, 49, 60], our scheme returns exactly one bounding-box per class per training image. This enables clean comparisons to previous work in terms of CorLoc on the training set, and keeps the human verification tasks simple (as we do not need to ask the annotators whether they see additional instances in an image). Note how at test time the detector is capable of localizing multiple objects of the same class in the same image (and this is captured in the mAP measure).

Compared methods. We compare our approach to the fully supervised alternative by training the same object detector (sec. 4.1) on the same training images, but with manual bounding-boxes (again, one bounding-box per class per image). On the other end of the supervision spectrum, we also compare to a modern MIL-based WSOL technique (sec 4.2) run on the same training images, but without human verification. Since that technique also forms the initialization step of our method, this comparison reveals how much farther we can go with human verification.

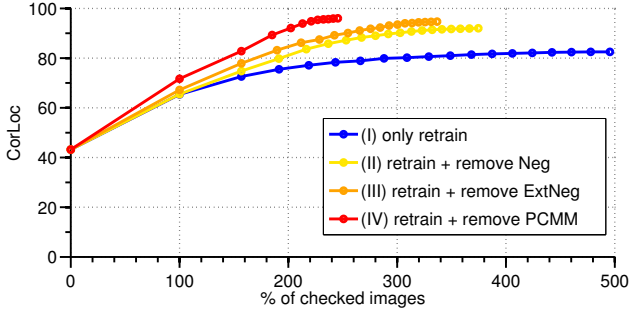


Figure 5: Trade-off between the number of verifications and CorLoc for the simulated verification case on PASCAL VOC 2007.

For MIL WSOL, the effort to draw bounding-boxes is zero. For fully supervised learning we take the actual annotation times for ILSVRC from [50]: they report 26 seconds for drawing one bounding-box without quality control, and 42 seconds with quality control. These timings are also representative for PASCAL VOC, since it is of comparable difficulty and its annotations are of comparable quality, as discussed in [40]. The bounding-boxes in both datasets are of high quality and precisely match the object extent.

5.2. Simulated verification

We first use simulated verification to determine how best to use the verification signal. We simulate human verification by using the available ground-truth bounding boxes. Note how these are *not* given to the learning algorithm, they are only used to derive the verification signals of sec. 3.1. Fig. 5 compares four ways to use the verification signal in terms of the trade-off between the number of verifications and CorLoc (sec. 3.3): **(I) only retrain** the object detector (using positively verified detections D_n^+); **(II) retrain + remove Neg**: for Yes/No verification, retrain and reduce the search space by eliminating one proposal for each negatively verified detection; **(III) retrain + remove ExtNeg**: for Yes/No verification, retrain and eliminate all proposals overlapping with a negatively verified detection; **(IV) retrain + remove PCMM**: for YPCMM verification, retrain and eliminate proposals according to the type of error.

As fig. 5 shows, even using verification just to retrain the object detector (I) drastically boosts CorLoc from the initial 43% (achieved by MIL WSOL) up to 82%. This requires checking each training image on average 4 times. Using the verification signal in the re-localisation step by reducing the search space (II–IV) helps to reach this CorLoc substantially faster (1.6–2 checks per image). Moreover, the final CorLoc is much higher when we reduce the search space. Removing negatively verified detections brings a considerable jump to 92% CorLoc (II); removing all proposals around negatively verified detections further increases CorLoc to 95% (III); while the Yes/Part/Container/Mixed/Missed strategy appears to be the most promising, achieving a near-perfect 96% CorLoc af-

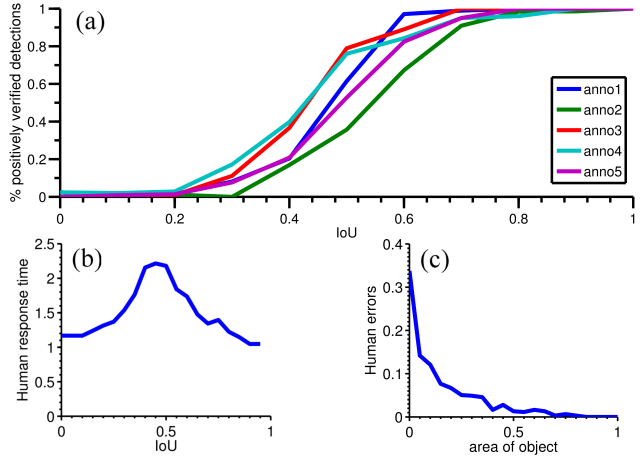


Figure 6: (a) Percentage of positively verified detections as a function of ground-truth IoU, for each annotator. (b) Average human response time as a function of ground-truth IoU. (c) Percentage of incorrectly verified detections as a function of object area (relative to image area).

ter checking each image only 2.5 times on average. These results show that feeding the verification signal into both the re-training and re-localisation steps quickly results in a large number of correctly located object instances.

Fig. 4 compares the search process of Yes/No and YPCMM verification strategies on a bird example. The second detection is diagnosed in YPCMM as a container. This focuses the search to that particular part of the image. In contrast, detections of the Yes/No case jump around before finding the detection. This shows that the search process of YPCMM is more targeted. However, in both cases the target object location is found rather quickly.

In conclusion, YPCMM is the most promising verification strategy, followed by Yes/No with removing all proposals overlapping with a negatively verified detection (III). Since Yes/No verification is intuitively easier and faster, we try both strategies in experiments with human annotators.

5.3. Human verification

Instructions and interface. For the Yes/No and YPCMM verification tasks, we used five annotators from our university who were given examples to learn about the IoU criterion. For both tasks we created a full-screen interface. All images of a target class were shown in sequence, with the current detection superimposed. For Yes/No verification, annotators were asked to press “1” for Yes and “0” for No. This simple task took on average 1.6 seconds per verification. For YPCMM verification, the annotators were asked to click on one of five on-screen buttons corresponding to Yes, Part, Container, Mixed and Missed. This more elaborate task took on average 2.4 seconds per verification.

Analysis of human verification. Fig. 6a reports the percentage of positively verified detections as a function of their IoU with the ground-truth bounding-box. We observe

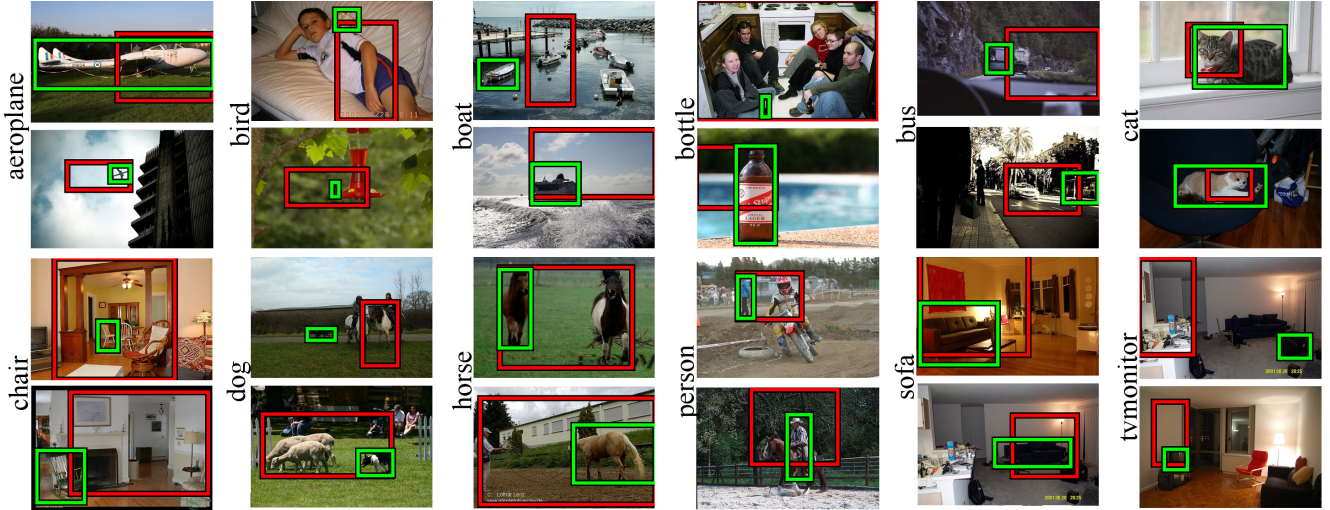


Figure 7: Examples of objects localized by using our proposed Yes/No human verification scheme on the trainval set of PASCAL VOC 2007 (sec. 5.3). For each example, we compare the final output of our scheme (green box) to the output of the reference multiple instance learning (MIL) weakly-supervised object localization approach (red box) (sec. 4.2).

that humans behave quite closely to the desired PASCAL criterion (i.e. $\text{IoU} > 0.5$) which we use in our simulations. All annotators behave identically on easy cases ($\text{IoU} < 0.25$, $\text{IoU} > 0.75$). On boundary cases ($\text{IoU} \approx 0.5$) we observe some annotator bias. For example, *anno2* tends to judge boundary cases as wrong detections, whereas *anno3* and *anno4* judge them more frequently as correct. Overall the percentage of incorrect Yes and No judgements are 14.8% and 8.5%, respectively. Therefore there is a slight bias towards Yes, i.e. humans tend to be slightly more lenient than the $\text{IoU} > 0.5$ criterion.

While the average human response time is 1.6 s for the Yes/No verification, the response time for verifying difficult detections ($\text{IoU} \approx 0.5$) is significantly higher (2.2 s, fig. 6b). This shows how the difficulty of the verification task is directly linked to the IoU of the detection, and is reflected in the time measurements. We also found that human verification errors strongly correlate with the area of objects: 48% of all errors are made when objects occupy less than 10% of the image area (fig. 6c).

Simulated vs. human annotators. We first compare simulated and human annotators by plotting CorLoc and mAP against actual annotation time (rather than as number of verifications as in Sec. 5.2). For simulated verification, we use average human annotation time as reported above. Fig. 8 shows the results on a log-scale. While human annotators are somewhat worse than simulations in terms of CorLoc on the training set, the mAP of the resulting object detectors on the test set are comparable. The diverging results for CorLoc and mAP is because human judgement errors are generally made on boundary cases with bounding-boxes that only approximately cover an object (fig. 6a). Using these cases either as positive or negative training examples, the object detector remains equally strong. To conclude, in terms of

training high quality object detectors, simulated human annotators reliably deliver similar results as actual annotators.

In sec. 5.2 we observed that YPCMM needs fewer verifications than Yes/No. However, in terms of total annotation time, the Yes/No task has the more favourable trade-off: Yes/No achieves 83% CorLoc and 45% mAP by taking 5.8 hours of annotation time, while YPCMM achieves 81% CorLoc and 45% mAP by taking 7.7 hours (fig. 8). Hence we conclude that the Yes/No task is preferable for human annotation, as it is easier and faster.

Weak supervision vs. verification. We now compare the reference MIL WSOL approach that we use to initialize our process (Sec. 4.2 and magenta diamond in fig. 8) to the final output of our Yes/No human verification scheme (solid orange line, fig. 8). While MIL WSOL achieves 43% CorLoc and 27% mAP, using human verification bring a massive jump in performance to 83% CorLoc and 45% mAP. Hence at a modest cost of 5.8 hours of annotation time we achieve substantial performance gains. Examples in fig. 7 show that our approach localizes objects more accurately and succeeds in more challenging conditions, e.g. when the object is very small and appears in a cluttered scene.

The state-of-the-art WSOL approaches perform as follows: Cinbis et al. [9] achieve 52.0% CorLoc and 30.2% mAP, Bilen et al. [5] achieve 43.7% CorLoc and 27.7% mAP, and Wang et al. [60] achieve 48.5% CorLoc and 31.6% mAP. Our method using human verification substantially outperforms all of them, reaching 83% CorLoc and 45% mAP (our method and [5, 9, 60] all use AlexNet). Hence at a modest extra annotation cost, we obtain many more correct object locations and train better detectors.

Full supervision vs. verification We now compare our Yes/No human verification scheme (solid orange line, fig. 8) to standard fully supervised learning with manual

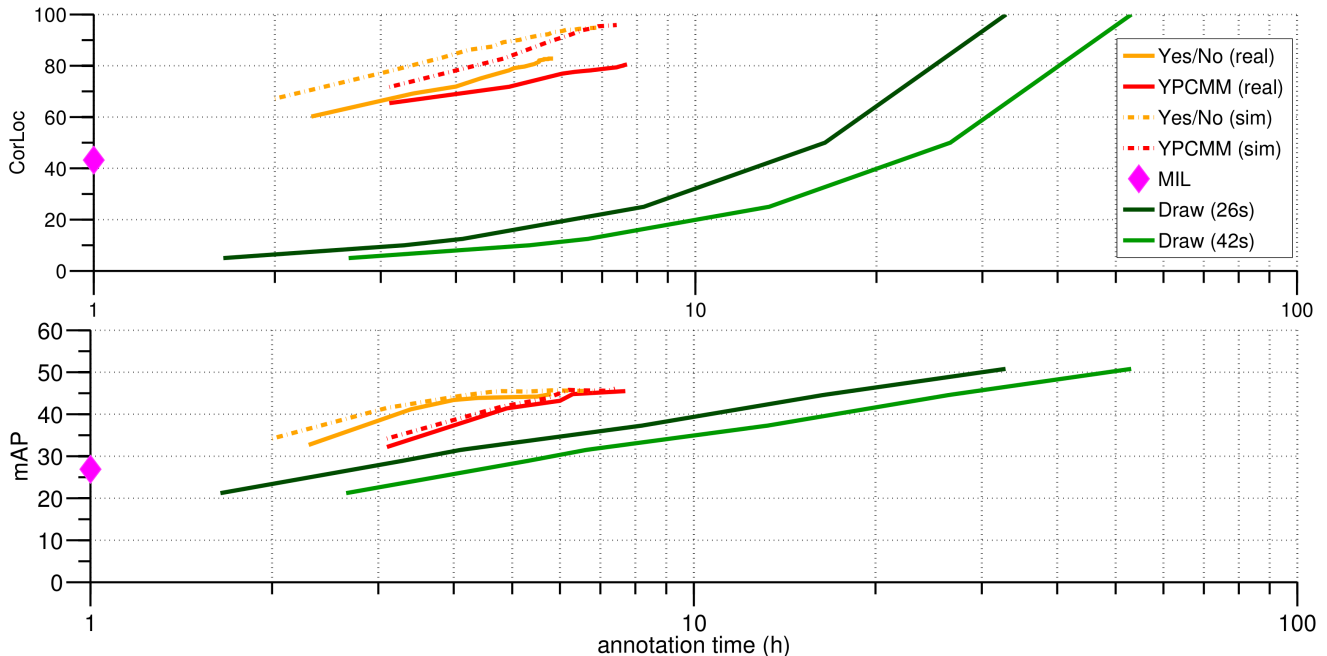


Figure 8: Evaluation on PASCAL VOC 2007: CorLoc and mAP against human annotation time in hours (log-scale). All orange and red curves are variations of our proposed scheme, with simulated (‘sim’) and real (‘real’) annotators. ‘Draw’ indicates learning from manual bounding-boxes (full supervision). ‘MIL’ indicates learning under weak supervision only, without human verification (sec. 4.2). The fully supervised approach needs a factor $6\times-9\times$ extra annotation time to obtain similar performance to our framework.

bounding-boxes (solid green lines). The object detectors learned by our scheme achieve 45% mAP, almost as good as the fully supervised ones (51% mAP). Importantly, fully supervised training needs 33 hours of annotation time (when assuming an optimistic 26 s per image), or even 53 hours (when assuming a more realistic 42 s per image). Our method instead requires only 5.8 hours, a reduction in human effort of a factor of $6\times-9\times$.

From a different perspective, when given the same human annotation time as our approach (5.8 hours), the fully supervised detector only achieves 33% mAP (at 26 s per bounding-box) or 30% mAP (at 42 s).

We conclude that by the use of just an inexpensive verification component, we can train strong object detectors at little cost. This is significant since it enables the cheap creation of high quality object detectors for a large variety of classes, bypassing the need for massive annotation efforts such as ImageNet [40].

5.4. Complete training set and VGG16

Our experiments are based on the Fast R-CNN detector [19]. In order to have a clean comparison between our verification-based scheme and the fully supervised results of [19], we re-ran our experiments using the complete trainval set of PASCAL VOC 2007 (i.e. 5011 images, table 1). Under full supervision, [19] reports 57% mAP based on AlexNet. Training Fast R-CNN from one bounding-box per class per image, results in 55% mAP, while our Yes/No human verification scheme gets to 50% mAP. Additionally, we

	reduced training set		complete training set	
	Yes/No	FS	Yes/No	FS
AlexNet	45%	51%	50%	55%
VGG16	55%	61%	58%	66%

Table 1: Comparison of mAP results between our Yes/No human verification scheme and full supervision (FS) using different training sets and different network architectures. ‘reduced training set’: excluding trainval images containing only difficult and truncated instances (3550 images); ‘complete training set’: all trainval images (5011).

experiment with VGG16 instead of AlexNet, with the same settings. Training with full supervision leads to 66% mAP, while our verification scheme delivers 58% mAP. Hence, on both CNN architectures our verification-based training scheme produces high quality detectors, achieving 90% of the mAP of their fully supervised counterparts.

5.5. Conclusions

We proposed a scheme for training object class detectors which uses a human verification step to improve the re-training and re-localisation steps common to most weakly supervised approaches. Experiments on PASCAL VOC 2007 show that our scheme produces detectors performing almost as good as those trained in a fully supervised setting, *without ever drawing any bounding-box*. As the verification task is very quick, our scheme reduces the total human annotation time by a factor of $6\times-9\times$.

Acknowledgement. This work was supported by the ERC Starting Grant “VisCul”.

References

- [1] B. Alexe, T. Deselaers, and V. Ferrari. What is an object? In *CVPR*, 2010. 3, 4, 5
- [2] Y. Aytar and A. Zisserman. Tabula rasa: Model transfer for object category detection. In *ICCV*, 2011. 3
- [3] T. Berg, A. Berg, J. Edwards, M. Mair, R. White, Y. Teh, E. Learned-Miller, and D. Forsyth. Names and Faces in the News. In *CVPR*, 2004. 3
- [4] H. Bilen, M. Pedersoli, and T. Tuytelaars. Weakly supervised object detection with posterior regularization. In *BMVC*, 2014. 1, 2, 3, 5
- [5] H. Bilen, M. Pedersoli, and T. Tuytelaars. Weakly supervised object detection with convex clustering. In *CVPR*, 2015. 1, 2, 3, 5, 7
- [6] A. Biswas and D. Parikh. Simultaneous active learning of classifiers & attributes via relative feedback. In *CVPR*, 2013. 2
- [7] S. Branson, C. Wah, F. Schroff, B. Babenko, P. Welinder, P. Perona, and S. Belongie. Visual recognition with humans in the loop. In *ECCV*, 2010. 2
- [8] R. Cinbis, J. Verbeek, and C. Schmid. Multi-fold mil training for weakly supervised object localization. In *CVPR*, 2014. 2, 3, 5
- [9] R. Cinbis, J. Verbeek, and C. Schmid. Weakly supervised object localization with multi-fold multiple instance learning. *arXiv:1503.00949*, 2015. 1, 5, 7
- [10] N. Dalal and B. Triggs. Histogram of Oriented Gradients for human detection. In *CVPR*, 2005. 1
- [11] J. Deng, J. Krause, and L. Fei-Fei. Fine-grained crowdsourcing for fine-grained recognition. In *CVPR*, 2013. 2
- [12] T. Deselaers, B. Alexe, and V. Ferrari. Localizing objects while learning their appearance. In *ECCV*, 2010. 1, 2, 3, 5
- [13] T. G. Dietterich, R. H. Lathrop, and T. Lozano-Perez. Solving the multiple instance problem with axis-parallel rectangles. *Artificial Intelligence*, 89(1-2):31–71, 1997. 2
- [14] P. Dollar and C. Zitnick. Edge boxes: Locating object proposals from edges. In *ECCV*, 2014. 3, 4, 5
- [15] P. Duygulu, K. Barnard, N. de Freitas, and D. Forsyth. Object recognition as machine translation: Learning a lexicon for a fixed image vocabulary. In *ECCV*, 2002. 3
- [16] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The PASCAL Visual Object Classes (VOC) Challenge. *IJCV*, 2010. 1, 3, 5
- [17] L. Fei-Fei, R. Fergus, and P. Perona. Learning generative visual models from few training examples: An incremental bayesian approach tested on 101 object categories. *cviu*, 2007. 3
- [18] P. Felzenszwalb, R. Girshick, D. McAllester, and D. Ramanan. Object detection with discriminatively trained part based models. *IEEE Trans. on PAMI*, 32(9), 2010. 1
- [19] R. Girshick. Fast R-CNN. In *ICCV*, 2015. 1, 4, 8
- [20] R. Girshick, J. Donahue, T. Darrell, and J. Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *CVPR*, 2014. 1, 2, 5
- [21] M. Guillaumin and V. Ferrari. Large-scale knowledge transfer for object localization in imagenet. In *CVPR*, 2012. 3, 5
- [22] A. Gupta and L. Davis. Beyond nouns: Exploiting prepositions and comparators for learning visual classifiers. In *ECCV*, 2008. 3
- [23] K. He, X. Zhang, S. Ren, and J. Sun. Spatial pyramid pooling in deep convolutional networks for visual recognition. In *ECCV*, 2014. 4
- [24] J. Hoffman, S. Guadarrama, E. Tzeng, R. Hu, and J. Donahue. LSDA: Large scale detection through adaptation. In *NIPS*, 2014. 3
- [25] A. J. Joshi, F. Porikli, and N. Papanikolopoulos. Multi-class active learning for image classification. In *CVPR*, 2009. 2
- [26] A. Kapoor, K. Grauman, R. Urtasun, and T. Darrell. Active learning with gaussian processes for object categorization. In *ICCV*, 2007. 2
- [27] A. Kovashka, S. Vijayanarasimhan, and K. Grauman. Actively selecting annotations among objects and attributes. In *ICCV*, 2011. 2
- [28] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*, 2012. 2, 4, 5
- [29] D. Kuettel, M. Guillaumin, and V. Ferrari. Segmentation Propagation in ImageNet. In *ECCV*, 2012. 3
- [30] S. Lad and D. Parikh. Interactively guiding semi-supervised clustering via attribute-based explanations. In *ECCV*, 2014. 2
- [31] C. Lampert, H. Nickisch, and S. Harmeling. Learning to detect unseen object classes by between-class attribute transfer. In *CVPR*, 2009. 3
- [32] C. Leistner, M. Godec, S. Schuster, A. Saffari, and H. Bischof. Improving classifiers with weakly-related videos. In *CVPR*, 2011. 3
- [33] J. Luo, B. Caputo, and V. Ferrari. Who’s doing what: Joint modeling of names and verbs for simultaneous face and pose annotation. In *NIPS*, 2009. 3
- [34] T. Malisiewicz, A. Gupta, and A. Efros. Ensemble of exemplar-svms for object detection and beyond. In *ICCV*, 2011. 1
- [35] D. P. Papadopoulos, A. D. F. Clarke, F. Keller, and V. Ferrari. Training object class detectors from eye tracking data. In *ECCV*, 2014. 3
- [36] D. Parikh and K. Grauman. Relative attributes. In *ICCV*, 2011. 2
- [37] A. Parkash and D. Parikh. Attributes for classifier feedback. In *ECCV*, 2012. 2
- [38] A. Prest, C. Leistner, J. Civera, C. Schmid, and V. Ferrari. Learning object class detectors from weakly annotated video. In *CVPR*, 2012. 3, 5
- [39] G.-J. Qi, X.-S. Hua, Y. Rui, J. Tang, and H.-J. Zhang. Two-dimensional active learning for image classification. In *CVPR*, 2008. 2
- [40] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. Berg, and L. Fei-Fei. Imagenet large scale visual recognition challenge. *IJCV*, 2015. 1, 2, 6, 8
- [41] O. Russakovsky, L.-J. Li, and L. Fei-Fei. Best of both worlds: human-machine collaboration for object annotation. In *CVPR*, 2015. 2

- [42] O. Russakovsky, Y. Lin, K. Yu, and L. Fei-Fei. Object-centric spatial pooling for image classification. In *ECCV*, 2012. 1, 2, 3, 5
- [43] N. Shapovalova, A. Vahdat, K. Cannons, T. Lan, and G. Mori. Similarity constrained latent support vector machine: An application to weakly supervised action classification. In *ECCV*, 2012. 5
- [44] Z. Shi, P. Siva, and T. Xiang. Transfer learning by ranking for weakly supervised object annotation. In *BMVC*, 2012. 2, 5
- [45] B. Siddiquie and A. Gupta. Beyond active noun tagging: Modeling contextual interactions for multi-class active learning. In *CVPR*, 2010. 2
- [46] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015. 4, 5
- [47] P. Siva and T. Xiang. Weakly supervised object detector learning with model drift detection. In *ICCV*, 2011. 1, 2, 3, 5
- [48] H. Song, R. Girshick, S. Jegelka, J. Mairal, Z. Harchaoui, and T. Darrell. On learning to localize objects with minimal supervision. In *ICML*, 2014. 1, 2, 3, 5
- [49] H. Song, Y. Lee, S. Jegelka, and T. Darrell. Weakly-supervised discovery of visual pattern configurations. In *NIPS*, 2014. 1, 2, 3, 5
- [50] H. Su, J. Deng, and L. Fei-Fei. Crowdsourcing annotations for visual object detection. In *AAAI Human Computation Workshop*, 2012. 1, 6
- [51] K. Tang, A. Joulin, L.-J. Li, and L. Fei-Fei. Co-localization in real-world images. In *CVPR*, 2014. 5
- [52] K. Tang, R. Sukthankar, J. Yagnik, and L. Fei-Fei. Discriminative segment annotation in weakly labeled video. In *CVPR*, 2013. 3
- [53] J. R. R. Uijlings, K. E. A. van de Sande, T. Gevers, and A. W. M. Smeulders. Selective search for object recognition. *IJCV*, 2013. 1, 3, 4
- [54] S. Vijayanarasimhan and K. Grauman. Multi-level active prediction of useful image annotations for recognition. In *NIPS*, 2008. 2
- [55] S. Vijayanarasimhan and K. Grauman. What's it going to cost you?: Predicting effort vs. informativeness for multi-label image annotations. In *CVPR*, 2009. 2
- [56] S. Vijayanarasimhan and K. Grauman. Large-scale live active learning: Training object detectors with crawled data and crowds. *IJCV*, 108(1-2):97–114, 2014. 1, 2
- [57] P. A. Viola, J. Platt, and C. Zhang. Multiple instance boosting for object detection. In *NIPS*, 2005. 1
- [58] C. Wah, S. Branson, P. Perona, and S. Belongie. Multiclass recognition and part localization with humans in the loop. In *ICCV*, 2011. 2
- [59] C. Wah, G. Van Horn, S. Branson, S. Maji, P. Perona, and S. Belongie. Similarity comparisons for interactive fine-grained categorization. In *CVPR*, 2014. 2
- [60] C. Wang, W. Ren, J. Zhang, K. Huang, and S. Maybank. Large-scale weakly supervised object localization via latent category learning. *IEEE Transactions on Image Processing*, 24(4):1371–1385, 2015. 1, 2, 3, 5, 7
- [61] L. Wang, G. Hua, R. Sukthankar, J. Xue, and J. Zheng. Video object discovery and co-segmentation with extremely weak supervision. In *ECCV*, 2014. 5
- [62] X. Wang, M. Yang, S. Zhu, and Y. Lin. Regionlets for generic object detection. In *ICCV*, pages 17–24. IEEE, 2013. 1
- [63] A. Yao, J. Gall, C. Leistner, and L. Van Gool. Interactive object detection. In *CVPR*, 2012. 1, 2